

# A program for the Bayesian Neural Network in the ROOT framework

Jiahang Zhong<sup>a,b</sup>, Run-Sheng Huang<sup>a</sup>, Shih-Chang Lee<sup>b,\*</sup>

<sup>a</sup>*School of Physics, Nanjing University, CN - Nanjing 210093, China*

<sup>b</sup>*Institute of Physics, Academia Sinica, TW - Taipei 11529, Taiwan*

---

## Abstract

We present a Bayesian Neural Network algorithm implemented in the TMVA package [1], within the ROOT framework [2]. Comparing to the conventional utilization of Neural Network as discriminator, this new implementation has more advantages as a non-parametric regression tool, particularly for fitting probabilities. It provides functionalities including cost function selection, complexity control and uncertainty estimation. An example of such application in High Energy Physics is shown. The algorithm is available with ROOT release later than 5.29.

*Keywords:* Bayesian Neural Network; TMVA; ROOT; Regression

---

## PROGRAM SUMMARY

*Manuscript Title:* A program for the Bayesian Neural Network in the ROOT framework

*Authors:* Jiahang Zhong, Run-Sheng Huang, Shih-Chang Lee

*Program Title:* TMVA-BNN

*Journal Reference:*

*Catalogue identifier:*

*Licensing provisions:* BSD

*Programming language:* C++

*Computer:* Any computer system or cluster with C++ compiler and UNIX-like operating system.

*Operating system:* Most UNIX/Linux systems. The application programs were thoroughly tested under Fedora and Scientific Linux CERN.

*Keywords:* Bayesian Neural Network, TMVA, ROOT, Regression.

*Classification:* 11.9

*External routines/libraries:* ROOT package version 5.29 or higher (<http://root.cern.ch>)

*Nature of problem:* Non-parametric fitting of multivariate distributions.

---

\*Corresponding author.

*E-mail address:* phsclee@phys.sinica.edu.tw

*Solution method:* An implementation of Neural Network following the Bayesian statistical interpretation. Uses Laplace approximation for the Bayesian marginalizations. Provides the functionalities of automatic complexity control and uncertainty estimation.

*Running time:* Time consumption for the training depends substantially on the size of input sample, the NN topology, the number of training iterations, etc. For the example in this manuscript, about 7 minutes was used on a PC/Linux with 2.0GHz processors.

## 1. Introduction

Neural Network (NN) has been considerably utilized in High Energy Physics in the past decade. In most applications, The NN was used as a discriminator to separate signal from backgrounds. It is a powerful tool to extract the features of target categories in the multivariate phase space, and project them into a scalar discriminator. However, such applications are often criticized for the reliance on the simulation or test-beam data as training samples, which may have distinct features comparing to the real data. On the other hand, the usage of NN as a non-parametric regression tool is much less exploited, and usually does not suffer from such concerns. Given sufficient complexity, even a single-hidden-layer NN can be seen as a universal approximator of any nonlinear multivariate function [3, 4]. Composed of simple nonlinear functions called “neurons”, such as hyperbolic tangent functions, an NN can achieve great complexity by connecting many neurons with variable weights  $\mathbf{w}$ , which serves as the free parameters of the model. Equation (1) shows the analytical form of a typical NN, with its structure shown in figure 1.

$$\begin{aligned}
 y_i^1 &= f^{(1)}(x_i) & f^{(1)}(h) &= ah + b \\
 y_j^2 &= f^{(2)}(w_{0j}^1 + \sum_i w_{ij}^1 y_i^1) & f^{(2)}(h) &= \tanh(h) \\
 y_1^3 &= f^{(3)}(w_{01}^2 + \sum_j w_{j1}^2 y_j^2) & f^{(3)}(h) &= h
 \end{aligned} \tag{1}$$

NN can approximate not only functions whose output values span real number space, but also those with confined output values. One such category particularly interesting is a probability. The output value can be confined within 0 and 1 by applying a sigmoid transformation to the output neuron, i.e. replacing  $f^{(3)}$  in equation (1) by

$$f^{(3)}(h) = \frac{1}{1 + e^{-h}} \tag{2}$$

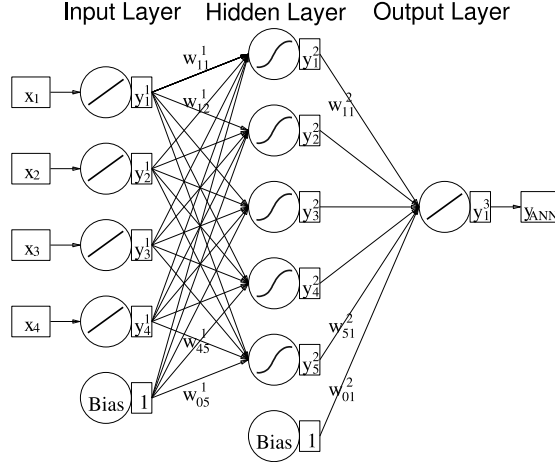


Figure 1: The structure of an NN with one hidden layer [1].

There are several advantages of NN comparing to the conventional representation of probabilities by histograms. First, NN can approximate the distributions in an unbinned manner, without the arbitrariness in the choice of binning and the subsequent loss of information. Second, NN is more practical for multivariate approximation without suffering the curse of dimensionality. Another great advantage of NN as a regression tool is that the correlations between the input variables can be well approximated.

Compared to the application as a “black-box” discriminator, NN application as a regression tool must follow a more explicit statistical interpretation. This is particularly important if the subsequent application requires detailed statistical information, such as limit setting for new physics searches. Here we present an NN algorithm following the Bayesian inference theory, implemented in the TMVA package [1]. This manuscript is organized as follows: section 2 is a brief review of the statistical interpretation for the training (fitting) and prediction procedure of the NN. Then the implementation of our Bayesian NN (BNN) algorithm is described in detail in section 3. Finally in section 4, an application of this BNN in High Energy Physics is demonstrated, where it is used to approximate the false identification rate of isolated muons.

## 2. Statistical Interpretation of Training and Prediction

As a generic model, NN normally has a huge number of degrees of freedom and incomprehensible complexity. On the other hand, the procedure to train an NN and makes predictions with it can be clearly interpreted by probability theory as Bayesian inference.

Given an observed sample  $\mathbf{D} = \{\mathbf{x}_i, t_i\}$ , with  $\mathbf{x}_i$  as the multiple input variables of entry  $i$ , and  $t_i$  as its observed output value, the training (fitting) of

NN can be viewed as a process to determine the probability of free parameters  $\mathbf{w}$ , based on  $\mathbf{D}$ . According to the Bayes theorem, this “posterior” probability comes from the combination of previous knowledge of  $\mathbf{w}$  (“prior”) and the compatibility of sample  $\mathbf{D}$  with the NN (“likelihood”)

$$P(\mathbf{w}|\mathbf{D}) \propto P(\mathbf{D}|\mathbf{w})P(\mathbf{w}) \quad (3)$$

The likelihood  $P(\mathbf{D}|\mathbf{w})$  in Bayesian language is closely related the “cost function” in machine learning language. The optimization of NN by minimizing the cost function is mostly equivalent to maximizing the likelihood. For example, the most commonly used sum of the square error (SSE) function can actually be translated as the negative logarithm of the Gaussian likelihood, as shown by equation (4).

$$\begin{aligned} \text{SSE} &= \sum_i (y(\mathbf{x}_i; \mathbf{w}) - t_i)^2 \\ &= -\log\left(\prod_i \exp(-(y(\mathbf{x}_i; \mathbf{w}) - t_i)^2)\right) \\ &\propto -\log(P(\mathbf{D}|\mathbf{w})) \end{aligned} \quad (4)$$

The prior probability  $P(\mathbf{w})$  is much less emphasized in classic usage of NN, which in fact often assumes a flat distribution. We will see later that a Bayesian regulator term can be added to the cost function, based on a simple prior knowledge. It is also worth mentioning that, although the probability distribution of  $\mathbf{w}$  can be obtained, only the most probable value is kept in classic usage.

Given a new input vector  $\mathbf{x}'$ , the prediction with NN can also be performed as a Bayesian inference. Although classic usage of NN only gives one single value of  $y'$  with the most probable  $\mathbf{w}$ , we should be able to predict a probability distribution of the output if we marginalize over the distribution of NN parameters  $P(\mathbf{w}|\mathbf{D})$  by equation (5).

$$P(y'|D, \mathbf{x}') = \int P(y'|\mathbf{x}', \mathbf{w})P(\mathbf{w}|D) d\mathbf{w} \quad (5)$$

### 3. Bayesian Implementations

#### 3.1. Cost Function

As shown in equation (4), the commonly used cost function, sum of the square error, can be interpreted as the negative logarithm of Gaussian likelihood function. This cost function is applicable for most regression applications, where a Gaussian distribution can be assumed for the observed values around their true values.

However, when the NN is used to approximate a probabilistic distribution, such an assumption of Gaussian likelihood may not be appropriate. As a probabilistic classifier, the NN’s output  $y$  is expected to approximate the probability of membership to one category, constrained between 0 and 1 by the sigmoid

function (equation (2)). And the observed value  $t$  for each entry in the input sample  $D$  is usually either 0 or 1, representing the fact whether the entry meets the condition, or belongs to the desired category. The distribution of observation  $t$  around probability  $y$  should then follow the Bernoulli distribution

$$P(\mathbf{D}|\mathbf{w}) = y^t(1 - y)^{1-t} \quad (6)$$

Correspondingly, the cost function for classification should take the form of the so-called “cross-entropy” function (CE), the sum of negative logarithm of Bernoulli distribution,

$$\text{CE} = \sum_i (-t_i \log y(\mathbf{x}_i; \mathbf{w}) - (1 - t_i) \log(1 - y(\mathbf{x}_i; \mathbf{w}))) \quad (7)$$

With the sigmoid transformation and CE cost function, the NN can approximate the probabilities with rigorous statistics interpretation. This can be chosen by the option of the MLP method [1] `EstimatorType=CE`.

### 3.2. Complexity Control

NN gains the capability of universal approximation by a huge number of degrees of freedom. A typical single-hidden-layer NN, even only for a few input variables, could have  $O(10^2)$  free parameters  $\mathbf{w}$ . A model with such great complexity may suffer from over-fitting. That is to say, it will approximate not only the desired connection between the input variables and the output value, but also the undesired fluctuations of the input sample. This is particularly an issue if the input sample has limited statistics. The predictivity of the model will be greatly deteriorated by exaggerated fluctuations.

The commonly employed solution against over-fitting is the so-called “cross validation” technique. A fraction of the input sample, normally half of the statistics, is taken away from the training process and used as a test set. Over-fitting is identified during the process of training, if the cost function of the test set starts to increase. One problem with this technique is the possibility that the cost function of test set may have merely hit a local minimum. Another problem is the reduction of training sample size. This is a non-trivial drawback for cases in which over-fitting may occur, which normally have trouble with statistics already.

In our implementation, another solution with regulators is adopted to avoid over-fitting. Although it is necessary to keep a large number of free parameters in order to make the model generic, the value of the parameters can be constrained to reduce unnecessary complexity. This can be expressed as a Bayesian prior knowledge about the model, assuming the value of the free parameters  $\mathbf{w}$  should be limited to the vicinity of zero. A Gaussian distribution centered at zero is used to represent such prior. Correspondingly, a “regulator” function can be obtained as the negative logarithm of this Gaussian prior for all NN parameters  $w_i$ .

$$\text{Reg} = -\log(P(\mathbf{w})) = \sum_i (\alpha_i w_i^2) \quad (8)$$

Adding this regulator term into the cost function actually gives the negative logarithm of posterior probability  $P(\mathbf{w}|D)$ . In the BNN, this summed value is minimized instead to obtain the optimal  $\mathbf{w}$ .

In equation (8), the hyper-parameters  $\alpha_i$  determine the range of  $\mathbf{w}$ , consequently reflecting the knowledge of required complexity of the model. The values of  $\alpha_i$  are not necessarily the same for each  $w_i$ . From a topological point of view, all the neurons within one hidden layer are computationally exchangeable. So their outgoing weights share the same hyper-parameter. This is not applicable to the input variables because they normally have different importance. Therefore, each group of weights originated from the same input neuron has its own  $\alpha_i$ . For the same reason, the bias neurons in each layer have independent hyper-parameters.

Although we can categorize the weights and associate them to different hyper-parameters, in most cases we do not actually possess *a priori* knowledge about the complexity needed, namely the values of  $\alpha_i$ . We implemented an iterative approach proposed by MacKay in 1992 [5], in which these hyper-parameters are estimated during the training of NN, by optimizing the “evidence” of the models:

$$P(D|\boldsymbol{\alpha}) = \int P(D|\mathbf{w}, \boldsymbol{\alpha}) P(\mathbf{w}|\boldsymbol{\alpha}) d\mathbf{w} \quad (9)$$

where the optimal  $\mathbf{w}$  is determined by minimizing the cost function and the integral is evaluated approximately for a given  $\boldsymbol{\alpha}$ . After estimating the optimal  $\boldsymbol{\alpha}$ , a new optimal  $\mathbf{w}$  is recalculated and the process is iterated, until desired convergence is reached.

This functionality can be activated by the MLP option `UseRegulator=true`, together with the BFGS training method.

### 3.3. Bayesian Prediction

In the classic usage of NN, the prediction upon new input  $\mathbf{x}'$  is the most probable value obtained with the most probable  $\mathbf{w}$ . Instead, in Bayesian data analysis, it is more essential to marginalize over all possible values of  $\mathbf{w}$ , and obtain the prediction as a probability distribution, as shown by equation (5). Such a prediction contains not only the most probable value, but also the uncertainty of the inference. The estimation of uncertainty is crucial, especially for extrapolated predictions in multivariate phase space.

Unfortunately, as a common difficulty for most Bayesian applications, the integration of equation (5) is generally non-trivial. Although the posterior of NN parameters  $P(\mathbf{w}|D)$  is calculable for any  $\mathbf{w}$ , the distribution for such “nuisance parameters” does not have a closed-form expression. In our implementation, an analytical distribution is used for the integration, which is the Laplace approximation of the posterior around the optimal value  $\mathbf{w}^{\text{MP}}$  [6], as shown in equation (10).

$$P(\mathbf{w}|D) \simeq P(\mathbf{w}^{\text{MP}}|D) \exp\left(-\frac{1}{2}\Delta\mathbf{w}^T \mathbf{A} \Delta\mathbf{w}\right) \quad (10)$$

$$\Delta\mathbf{w} = \mathbf{w} - \mathbf{w}^{\text{MP}}$$

$\mathbf{A}$  represents the Hessian matrix of the cost function, namely the negative logarithm of the posterior

$$\mathbf{A} = -\nabla\nabla \log P(\mathbf{w}|D)|_{\mathbf{w}^{\text{MP}}}. \quad (11)$$

For feed-forward NN either with a linear output neuron and SSE cost function, or with a sigmoid output neuron and CE cost function, its Hessian matrix  $\mathbf{A}$  can be approximated and written consistently as

$$\mathbf{A} \simeq \sum_{\mathbf{x}_i} f' \mathbf{g}^T \mathbf{g} \quad (12)$$

$f' = dy/dh$  and  $\mathbf{g} = \nabla h$ .  $h$  and  $y$  are the values before and after output neuron transformation.

Furthermore, a linear dependence of  $h$  over weights  $\mathbf{w}$  can be approximated as well.

$$h(\mathbf{x}'; \mathbf{w}) \simeq h(\mathbf{x}'; \mathbf{w}^{\text{MP}}) + \mathbf{g} \cdot \Delta \mathbf{w} \quad (13)$$

Therefore, the distribution of  $h$  can be calculated analytically as

$$\begin{aligned} P(h|D, \mathbf{x}') &= \int h(\mathbf{x}'; \mathbf{w}) P(\mathbf{w}|D) d\mathbf{w} \\ &\simeq \mathcal{N}(h(\mathbf{x}'; \mathbf{w}^{\text{MP}}), \mathbf{g}^T \mathbf{A}^{-1} \mathbf{g}) \end{aligned} \quad (14)$$

It is a Gaussian distribution with the mean value of the classic NN prediction. In addition, each prediction can give an associated uncertainty, which originates from the uncertainty of the NN parameters determination. For probability fitting, the probability density function (p.d.f.) of the output  $y$  is a sigmoid-transformed Gaussian distribution of  $h$  (equation 2). As a non-linear transformation, the sigmoid function will convert the symmetric error bar of  $h$  into an asymmetric one for  $y$ .

Besides the Laplace approximation used in our implementation, the integration can also be solved numerically by Markov Chain Monte Carlo (MCMC) [7]. Compared to MCMC approach, the analytical approximation is easier for both training and prediction. The estimation of the Hessian matrix during the training process can be activated by MLP option `CalculateErrors=true`. By configuring the Reader option `Error=true`, the Hessian matrix will be loaded for prediction. And the asymmetric uncertainties can be evaluated by Reader functions `GetMVAErrorUpper()` and `GetMVAErrorLower()`.

#### 4. Application in High Energy Physics

To demonstrate the practical usage of the BNN, we will show an example in High Energy Physics: measurement of the false identification rate of isolated muons. The test data, job control scripts and instructions can be obtained from the CPC library.

Muons are an important electroweak signature in collider physics. According to their sources, they can be categorized by the so-called isolation condition, i.e. adjacent particle flow. Those “non-isolated” muons, which have accompanying particles flying in a similar direction, mostly come from semi-leptonic decay of heavy flavor quarks ( $b, c$ ). The other “isolated” muons are more likely from electro-weak processes such as  $W$  boson production, and therefore taken as signals in many physics topics. To get rid of the former category from a mixed sample, an isolation cut is often imposed on the sum of transverse momenta, for all visible particles inside a cone around the muon. Here we choose a typical configuration, limiting transverse particle flow within  $\Delta R = \sqrt{(\Delta\eta)^2 + (\Delta\phi)^2} \leq 0.2$  to be less than 15% of the muon’s transverse momentum.

Due to event-specific kinematics and detector response, there is always a considerable fraction of muons from “non-isolated sources” being falsely identified as isolated muons. Such muons will contaminate the signal sample when a final state with isolated muons is expected. It is non-trivial to estimate this background accurately, for both new physics searches and precise measurements. Rather than relying on simulation, it is highly desirable to measure the false identification rate  $f = P(\text{isol})$  directly from collision data.

In the following, we demonstrate how BNN can be used for such false rate measurement, using Monte-Carlo simulation of heavy-flavor quark production ( $b\bar{b}, c\bar{c}$ ) in the proton-proton collisions. The samples are produced by the Pythia generator [8], with 7 TeV center-of-mass energy as at LHC [9]. Muon pseudo-rapidity acceptance is assumed to be within 2.5, with the threshold of transverse momentum as above 10 GeV. The fiducial efficiency of the detector is assumed to be 100% for simplicity.

Two samples are generated for comparison. The first sample requires single muon within the detector acceptance, while the second sample requires two such muons with the same charge. The former final state is dominated by the non-isolated sources, and therefore ideal for the measurement of fake rate. The latter final state instead is a typical channel in which new physics may manifest, and requires accurate estimation of the background contribution. With collision data, we need to measure the fake rate from the single muon events, then apply it to those same-charge di-muon events to estimate the background contribution [10]. With the MC samples mentioned above, we can study how to make the fake rate measurement compatible between these two channels.

In a first attempt to measure the average false identification rates, we obtain a quite different values,  $33.3(0.1)\pm\%$  from the single muon sample and  $19.8(0.2)\%$  from the same-charge di-muon sample. This incompatibility is due to the fact that the probability of false identification,  $P(\text{isol}|\mathbf{x})$ , is not constant over kinematic variables  $\mathbf{x}$ . Two examples of such kinematic variables are

- $p_T$ : the transverse momentum of the muon itself.
- $H_T$ : the scalar sum of transverse momenta for all visible particles in the event, and the measurable energy imbalance.

The kinematic distributions in these two samples are quite distinct, as can be



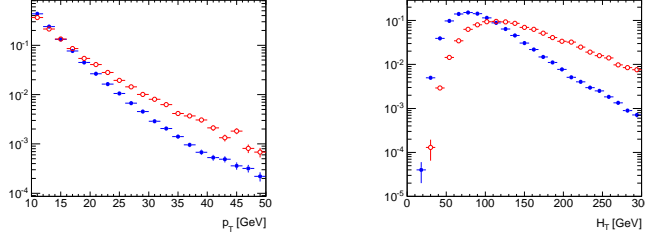


Figure 2: Probability density distributions over  $p_T$  (left) and  $H_T$  (right), for muons in the single muon control sample (solid) and di-muon signal sample (circle).

seen in figure 2. As a result, the observed average rates, marginalized by equation (15), turned out to be incompatible. In order to make a correct prediction in the signal region, it is important to measure the probabilities  $P(\text{isol} | \mathbf{x})$  rather than the marginalized rate.

$$\langle f \rangle = \int P(\text{isol} | \mathbf{x}) P(\mathbf{x}) d\mathbf{x} \quad (15)$$

Our BNN implementation is used to perform this measurement of  $P(\text{isol} | \mathbf{x})$ . About 50,000 muons in the single muon QCD events is used as input sample  $\mathbf{D}$ . Their corresponding kinematic variables  $p_T$  and  $H_T$  are declared as the input  $\mathbf{x}_i$  of the NN. And a target value  $t_i$  as 1 or 0 is assigned, depending on whether the muon passed the isolation criteria. The NN is constructed with one hidden layer of ten neurons, and a sigmoid-transformed output neuron. As a probability fit, the training is configured to use the cross-entropy cost function of equation (7). In addition, the regulator mechanism is activated to prevent over-fitting.

The fitted distribution  $P(\text{isol} | p_T, H_T)$  can be visualized in figure 3. The correlation between the two variables is clearly fitted, with a smooth extrapolation into the peripheral phase space region.

We then test this 2D false rate function with muons in the same-charge di-muon sample. Using the trained BNN, we can predict the probability of passing isolation criteria  $P(\text{isol} | p_T, H_T)$  for each muon. Marginalizing the predicted probabilities over all the muons in this sample (30764 muons), we can predict that 6077 of them will pass the isolation cut, corresponding to an average false rate of 19.8%. This is very close to the actual number of passed muon, 6093, equivalent to an average false rate of 19.8(0.2)%.

As described in section 3.3, the BNN can also calculate the uncertainty associated to each prediction, based on the uncertainty on the determination of the free parameters  $\mathbf{w}$ . It reflects the statistical property of the training sample. To test this estimation, another 100 NNs are trained with the same network topology. But the input samples, as well as the random seeds for initialization, are different in each training. For every muon in the same-charge di-muon sample, we use the BNN to predict its probability of pass  $P_{\text{BNN}}$ , as well as the associated asymmetric error bars, denoted as  $\sigma_{\text{BNN}}^+$  and  $\sigma_{\text{BNN}}^-$ . As a comparison, we also

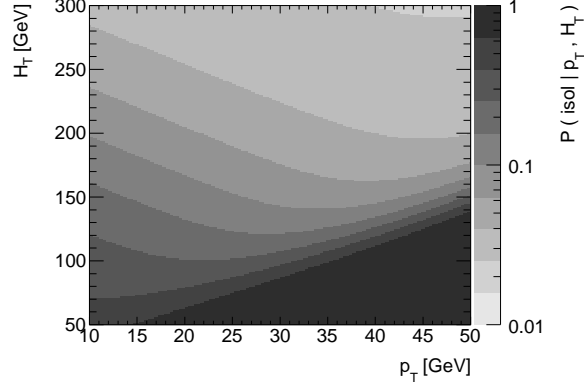


Figure 3: 2-dimensional false isolation probabilities, fitted by the BNN.

use the “batch” NNs to make 100 predictions, and calculate their mean value  $P_{\text{Batch}}$  and standard deviation  $\sigma_{\text{Batch}}$ . Comparing the ratio between  $\sigma_{\text{BNN}}$  and  $\sigma_{\text{Batch}}$  (figure 4(a)) for all the muons, we can see that the BNN uncertainty is generally consistent with the standard deviation of the batch predictions.

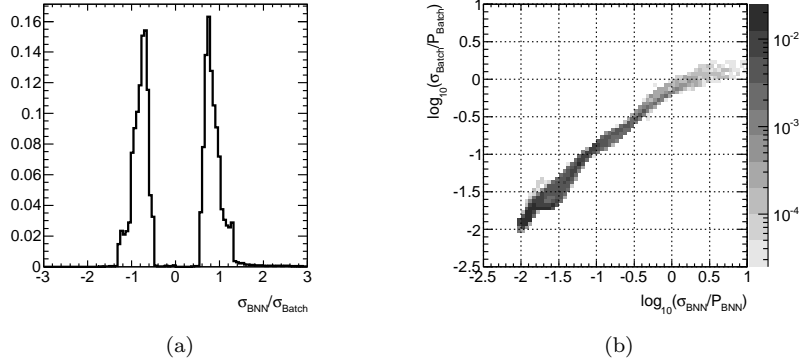


Figure 4: (a) The distribution of the ratios between  $\sigma_{\text{BNN}}$  and  $\sigma_{\text{Batch}}$ , for all muons in the same-charge di-muon events. Positive value stands for  $\sigma_{\text{BNN}}^+ / \sigma_{\text{Batch}}$  and negative value stands for  $\sigma_{\text{BNN}}^- / \sigma_{\text{Batch}}$ . (b) The correlation between  $\log_{10}(\sigma_{\text{BNN}} / P_{\text{BNN}})$  and  $\log_{10}(\sigma_{\text{Batch}} / P_{\text{Batch}})$ .  $\sigma_{\text{BNN}}$  is the average of  $\sigma_{\text{BNN}}^+$  and  $\sigma_{\text{BNN}}^-$ .

To further understand the precision of the BNN uncertainty estimation, the correlation between  $\log_{10}(\sigma_{\text{BNN}} / P_{\text{BNN}})$  and  $\log_{10}(\sigma_{\text{Batch}} / P_{\text{Batch}})$  for all the muons are plotted in figure 4(b). Good consistency can be observed when the relative uncertainties are less than  $\sim 30\%$ , which is the case for a large fraction of the entries. For predictions with large relative uncertainty, BNN tends to over-estimate the uncertainty, as the approximation applied in the Bayesian

marginalization becomes less accurate.

It is worthwhile to notice that the uncertainty estimated by BNN only accounts for the confidence in the determination of BNN parameters, evaluated based on the training sample used. The difference between the prediction and observation also involves the statistical fluctuations of the prediction sample, as well as the systematic uncertainties in the application, such as sample selections, the choice of parameterized variables and their measurement uncertainties.

Furthermore, only two kinematic variables were considered in this truth-level study, and consistent false rate estimation has already been observed. In reality, there could be additional factors which considerably affect the false rate, due to detector effects and collision configuration. Fortunately, the measurement can be easily extended to a higher dimensional phase space with BNN as the fitting tool.

## 5. Conclusion

In this manuscript we presented a BNN algorithm which can be used as an unbinned fitting tool, which is particularly interesting for fitting probabilities. The Bayesian implementation also provides functionalities such as controlling unnecessary complexity, and uncertainty estimation.

The demonstration with a HEP use case clearly showed the capability of BNN as an unbinned regression tool, especially if several input variables with correlation are involved. This technique has already been used to analyze the data collected by LHC in 2010 [10]. It has a very promising future for further applications, particularly for higher dimensional problems.

## Acknowledgements

We thank Andreas Hoecker, Joerg Stelzer, Peter Speckmayer, Jan Therhaag, Eckhard von Toerne and Helge Voss for their help in implementing this program into TMVA. We thank Song-Ming Wang and Zhili Weng for helpful discussions.

Jiahang Zhong and Shih-Chang Lee are partially supported by the National Science Council, Taiwan under the contract number NSC99-2119-M-001-015.

## References

- [1] A. Hoecker, P. Speckmayer, J. Stelzer, J. Therhaag, E. von Toerne, H. Voss, TMVA: Toolkit for Multivariate Data Analysis, PoS ACAT (2007) 040. [arXiv:physics/0703039](#).
- [2] R. Brun, F. Rademakers, Root – an object oriented data analysis framework, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment 389 (1-2) (1997) 81 – 86, new Computing Techniques in Physics Research V.

- [3] G. Cybenko, Approximation by Superpositions of a Sigmoidal function, *Mathematics of Control, Signals and Systems* 2 (1989) 303–314.
- [4] K. Hornik, M. Stinchcombe, H. White, Multilayer feedforward networks are universal approximators, *Neural Networks* 2 (5) (1989) 359 – 366.
- [5] D. J. C. MacKay, Bayesian interpolation, *Neural Computation* 4 (3) (1992) 415–447. doi:10.1162/neco.1992.4.3.415.
- [6] D. J. C. MacKay, Probable networks and plausible predictions - a review of practical bayesian methods for supervised neural networks, *Network: Computation in Neural Systems* 6 (3) (1995) 469–505.
- [7] R. M. Neal, Bayesian learning via stochastic dynamics, in: *Advances in Neural Information Processing Systems 5*, [NIPS Conference], Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1993, pp. 475–482.
- [8] T. Sjostrand, S. Mrenna, P. Skands, Pythia 6.4 physics and manual, *Journal of High Energy Physics* 2006 (05) (2006) 026.
- [9] L. Evans, P. Bryant, Lhc machine, *Journal of Instrumentation* 3 (08) (2008) S08001.
- [10] M. Aliev, *et al*, Inclusive search for exotic same-sign dilepton signatures at atlas, Tech. Rep. ATL-COM-PHYS-2011-107, CERN, Geneva (Feb 2011).